

# ENERGY STABLE FLUX FORMULAS FOR THE DISCONTINUOUS GALERKIN DISCRETIZATION OF FIRST ORDER NONLINEAR CONSERVATION LAWS

TIMOTHY BARTH\* AND PIERRE CHARRIER†

**Abstract.** We consider the discontinuous Galerkin (DG) finite element discretization of first order systems of conservation laws derivable as moments of the kinetic Boltzmann equation. This includes well known conservation law systems such as the Euler equations of gasdynamics. For the class of first order nonlinear conservation laws equipped with an entropy extension, an energy analysis of the DG method for the Cauchy initial value problem is developed. Using this DG energy analysis, several new variants of existing numerical flux functions are derived and shown to be energy stable.

**Key words.** Nonlinear Conservation Laws, Boltzmann Equation, Entropy Symmetrization, Discontinuous Galerkin Finite Element Method

**AMS subject classifications.** 35L02, 65M02, 65K02, 76N02

**1. Introduction.** Discontinuous Galerkin (DG) finite element methods for first order hyperbolic equations were introduced in the early works of Reed and Hill [19] and Johnson and Pitkäranta [15] with application to nonlinear conservation law systems by Cockburn *et al.* [4, 3, 5]. Fundamental to DG methods is the use of approximation spaces that are devoid of interelement continuity in both space and time. The multi-valued representation of the solution at interelement boundaries makes the evaluation of conservation law fluxes ambiguous thus necessitating the introduction of a *numerical flux function*,  $\mathbf{h}(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n})$ , a vector function of two (or more) solution states and a geometric normal at interelement boundaries. The needed numerical flux function can have design origins from exact or approximate solutions of the Riemann problem of gasdynamics [10, 20, 14]. Alternatively, the numerical flux can be designed from a nonlinear energy analysis of the DG method for first order nonlinear conservation laws equipped with a convex entropy extension, see Barth [2, 1]. This latter energy technique is used in the present analysis. Using the notation introduced in later consideration of the Cauchy initial-value problem, one obtains from this analysis the following *exact* energy balance equation for the DG finite element method for a spatial domain  $\Omega$  integrated over  $N$  time slabs

$$(1.1) \quad \frac{1}{2} \sum_{n=0}^{N-1} \left( \underbrace{\| [\mathbf{v}]_{t_n^+}^n \|_{\tilde{A}_0, \Omega}^2}_{\text{energy removal via discontinuous in time function representation}} + \underbrace{\sum_{e \in \mathcal{E}} \langle [\mathbf{v}]_{x^+}^n \rangle_{[\tilde{A}]_e}^2}_{\text{energy removal via discontinuous in space function representation}} \right) + \underbrace{\int_{\Omega} U(t_N^-) dx}_{\text{final energy}} = \underbrace{\int_{\Omega} U(t_0^-) dx}_{\text{initial energy}} .$$

This exact energy balance is derived using either of two different baseline numerical flux functions:

- **Symmetric Mean-Value (SMV) Flux**

$$\mathbf{h}_{\text{SMV}}(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) = \frac{1}{2}(\mathbf{f}(\mathbf{v}_-) + \mathbf{f}(\mathbf{v}_+)) - \frac{1}{2} \mathbf{h}_{\text{SMV}}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n})$$

with

$$\mathbf{h}_{\text{SMV}}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) \equiv \int_0^1 |\tilde{A}(\bar{\mathbf{v}}(\theta); \mathbf{n})|_{\tilde{A}_0} d\theta [\mathbf{v}]_-^+ .$$

- **Kinetic Symmetric Mean-Value (KSMV) Flux**

$$\mathbf{h}_{\text{KSMV}}(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) = \frac{1}{2}(\mathbf{f}(\mathbf{v}_-) + \mathbf{f}(\mathbf{v}_+)) - \frac{1}{2} \mathbf{h}_{\text{KSMV}}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n})$$

\*NASA Ames Research Center, Information Sciences Directorate, Moffett Field, California, 94035-1000 USA  
(barth@nas.nasa.gov)

†Mathématiques Appliquées de Bordeaux Université Bordeaux I, 351 Cours de la Libération, 33 405 Talence cedex, France  
(Pierre.Charrier@math.u-bordeaux.fr)

with

$$\mathbf{h}_{\text{KSMV}}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) \equiv \int_0^1 \langle |v \cdot \mathbf{n}| \mathbf{m} \otimes \mathbf{m} \exp(\bar{\mathbf{v}}(\theta) \cdot \mathbf{m}(v)) \rangle d\theta [\mathbf{v}]_-^+$$

where  $\langle \cdot \rangle$  denotes an integration in velocity-internal energy phase space and  $\mathbf{m}(v)$  the vector of moments as discussed later. Observe that the nonlinear energy balance (1.1) formally bounds the final solution in terms of initial data. The discontinuous function space leads to energy removal in space and time proportional to the matrix modulated square of solution jumps across the respective space and time interfaces.

In general, the numerical flux functions given here are too complex to permit the calculation of the needed path integrations in closed form. Our strategy in this paper is to first develop the framework and theoretical results given above. We then shift to our main objective, the development of approximate numerical flux functions,  $\mathbf{h}_{\text{approx}}(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n})$ , that avoid these complicated path integrations without compromising our ability to rigorously prove nonlinear stability. This is accomplished by requiring that the approximate numerical flux formulas are more energy dissipative than the theoretically derived fluxes given above. This task can be reduced to the satisfaction of either of two algebraic sufficient conditions (derived later)

$$[\mathbf{v}]_{x-}^{z+} \cdot \mathbf{h}_{\text{SMV}}^d \leq [\mathbf{v}]_{x-}^{z+} \cdot \mathbf{h}_{\text{approx}}^d \quad \text{or} \quad [\mathbf{v}]_{x-}^{z+} \cdot \mathbf{h}_{\text{KSMV}}^d \leq [\mathbf{v}]_{x-}^{z+} \cdot \mathbf{h}_{\text{approx}}^d .$$

We then construct a number of approximate flux functions based on this strategy.

**2. Background.** Consider the Cauchy initial value problem for a system of  $m$  coupled first-order differential equations in  $d$  space coordinates and time which represents a conservation law process. Let  $\mathbf{u}(x, t) : \mathbb{R}^d \times \mathbb{R}^+ \mapsto \mathbb{R}^m$  denote the dependent solution variables and  $\mathbf{f}(\mathbf{u}) : \mathbb{R}^m \mapsto \mathbb{R}^{m \times d}$  the flux vector. The prototype Cauchy problem is then given by

$$(2.1) \quad \begin{cases} \mathbf{u}_{,t} + \mathbf{f}_{,x_i}^i = 0 \\ \mathbf{u}(x, 0) = \mathbf{u}_0(x) \end{cases}$$

with implied summation on the index  $i$ . Additionally, the system is assumed to possess an scalar entropy extension. Let  $U(\mathbf{u}) : \mathbb{R}^m \mapsto \mathbb{R}$  denote an entropy function and  $F(\mathbf{u}) : \mathbb{R}^m \mapsto \mathbb{R}^d$  the entropy flux such that in addition to (2.1) the following inequality holds

$$(2.2) \quad U_{,t} + F_{,x_i}^i \leq 0$$

with equality for smooth solutions. In symmetrization theory for first-order conservation laws [11, 17, 12], one seeks a mapping  $\mathbf{u}(\mathbf{v}) : \mathbb{R}^m \mapsto \mathbb{R}^m$  applied to (2.1) so that when transformed

$$(2.3) \quad \mathbf{u}_{,\mathbf{v}} \mathbf{v}_{,t} + \mathbf{f}_{,\mathbf{v}}^i \mathbf{v}_{,x_i} = 0$$

the matrix  $\mathbf{u}_{,\mathbf{v}}$  is symmetric positive definite (SPD) and the matrices  $\mathbf{f}_{,\mathbf{v}}^i$  are symmetric. Clearly, if functions  $\mathcal{U}(\mathbf{v}) : \mathbb{R}^m \mapsto \mathbb{R}$  and  $\mathcal{F}^i(\mathbf{v}) : \mathbb{R}^m \mapsto \mathbb{R}$  can be found so that

$$(2.4) \quad \mathbf{u}^T = \mathcal{U}_{,\mathbf{v}}, \quad (\mathbf{f}^i)^T = \mathcal{F}_{,\mathbf{v}}^i$$

then the matrices

$$(2.5) \quad \mathbf{u}_{,\mathbf{v}} = \mathcal{U}_{,\mathbf{v},\mathbf{v}}, \quad \mathbf{f}_{,\mathbf{v}}^i = \mathcal{F}_{,\mathbf{v},\mathbf{v}}^i$$

are symmetric. Further, we shall require that  $\mathcal{U}(\mathbf{v})$  be a differentiable convex function such that

$$(2.6) \quad \lim_{\mathbf{v} \rightarrow \infty} \frac{\mathcal{U}(\mathbf{v})}{|\mathbf{v}|} = +\infty$$

so that  $U(\mathbf{u})$  can be interpreted as a Legendre transform of  $\mathcal{U}(\mathbf{v})$

$$(2.7) \quad U(\mathbf{u}) = \sup_{\mathbf{v}} \{ \mathbf{v} \cdot \mathbf{u} - \mathcal{U}(\mathbf{v}) \} .$$

From (2.6), it follows that  $\exists \mathbf{v}^* \in \mathbb{R}^m$  such that  $\mathbf{v} \cdot \mathbf{u} - \mathcal{U}(\mathbf{v})$  achieves a maximum at  $\mathbf{v}^*$

$$(2.8) \quad U(\mathbf{u}) = \mathbf{v}^* \cdot \mathbf{u} - \mathcal{U}(\mathbf{v}^*) .$$

At this maximum  $\mathbf{u} = \mathcal{U}_{,\mathbf{v}}(\mathbf{v}^*)$  which can be locally inverted to the form  $\mathbf{v}^* = \mathbf{v}(\mathbf{u})$ . Elimination of  $\mathbf{v}^*$  in (2.8) yields the simplified duality relationship

$$(2.9) \quad U(\mathbf{u}) = \mathbf{v}(\mathbf{u}) \cdot \mathbf{u} - \mathcal{U}(\mathbf{v}(\mathbf{u})) .$$

Differentiation of this expression

$$(2.10) \quad U_{,\mathbf{u}} = \mathbf{v}^T + \mathbf{u}^T \mathbf{v}_{,\mathbf{u}} - \mathcal{U}_{,\mathbf{v}} \mathbf{v}_{,\mathbf{u}} = \mathbf{v}^T$$

gives an explicit formula for the entropy variables  $\mathbf{v}$  in terms of derivatives of the entropy function  $U(\mathbf{u})$

$$(2.11) \quad \mathbf{v}^T = U_{,\mathbf{u}} .$$

Using the mapping relation  $\mathbf{v}(\mathbf{u})$ , a duality pairing for entropy flux components is defined

$$(2.12) \quad F^i(\mathbf{u}) = \mathbf{v}(\mathbf{u}) \cdot \mathbf{f}^i(\mathbf{u}) - \mathcal{F}^i(\mathbf{v}(\mathbf{u})) .$$

Differentiation then yields the flux relation

$$(2.13) \quad F_{,\mathbf{u}}^i = \mathbf{v}^T \mathbf{f}_{,\mathbf{u}}^i + (\mathbf{f}^i)^T \mathbf{v}_{,\mathbf{u}} - \mathcal{F}_{,\mathbf{v}}^i \mathbf{v}_{,\mathbf{u}} = \mathbf{v}^T \mathbf{f}_{,\mathbf{u}}^i$$

and the fundamental relationship for smooth solutions

$$(2.14) \quad \mathbf{v} \cdot (\mathbf{u}_{,t} + \mathbf{f}_{,x_i}^i) = U_{,t} + F_{,x_i}^i = 0$$

which is exploited in nonlinear energy analysis.

Note that convexity of  $\mathcal{U}(\mathbf{v})$  implies positive definiteness of  $\mathbf{u}_{,\mathbf{v}}$  and hyperbolicity of (2.1) [8, 17], viz., that the linear combination  $\mathbf{f}_{,\mathbf{u}}(\mathbf{n}) = n_i \mathbf{f}_{,\mathbf{u}}^i$  has real eigenvalues and a complete set of real-valued eigenvectors for all nonzero  $\mathbf{n} \in \mathbb{R}^d$ . This result follows immediately from the identity

$$(\mathbf{u}_{,\mathbf{v}})^{-1/2} \mathbf{f}_{,\mathbf{u}}(\mathbf{n})(\mathbf{u}_{,\mathbf{v}})^{1/2} = \underbrace{(\mathbf{u}_{,\mathbf{v}})^{-1/2} \mathbf{f}_{,\mathbf{v}}(\mathbf{n})(\mathbf{u}_{,\mathbf{v}})^{-1/2}}_{\text{symm}}$$

which shows that  $\mathbf{f}_{,\mathbf{u}}(\mathbf{n})$  is similar to a symmetric matrix.

**2.1. Kinetic Boltzmann Entropies.** Consider the particular case of moment systems derived from the kinetic Boltzmann equation with Levermore's closure [16]. Boltzmann's equation is given by

$$(2.15) \quad f(x, v, t)_{,t} + v \cdot \nabla_x f(x, v, t) = C(f)(x, v, t) ,$$

with  $f(x, v, t)$  a nonnegative density function,  $v \in \mathbb{R}^d$  the velocity, and  $C(f) : \mathbb{R} \mapsto \mathbb{R}$  the collision operator. Moment systems are obtained by integrating in velocity space the Boltzmann equation over a vector  $\mathbf{m}(\mathbf{v})$  of linearly independant polynomials in velocity,

$$(2.16) \quad \langle \mathbf{m} f \rangle_{,t} + \langle v_i \mathbf{m} f \rangle_{,x_i} = \langle \mathbf{m} C(f) \rangle ,$$

where  $\langle \psi \rangle$  denotes the integral of a measurable function  $\psi$  over velocity space. Without further assumption, the fluxes  $\langle v_i \mathbf{m} f \rangle$  cannot be expressed as functions of  $\mathbf{u} = \langle \mathbf{m} f \rangle$ . The closure of the system is performed by assuming that the distribution function  $f$  has a prescribed form  $f_B = f_B(\mathbf{u})$  given by the *minimum entropy principle*

$$(2.17) \quad H[f_B] = \min \{ H[g] \mid \langle g \mathbf{m} \rangle = \mathbf{u} \} ,$$

where  $H[g] = \langle g \ln g \rangle$  is Boltzmann's celebrated  $H$ -function. Since  $H$  is a convex function the minimization problem (2.17) is formally equivalent to

$$(2.18) \quad f_B = \exp(\mathbf{v} \cdot \mathbf{m}) ,$$

where  $\mathbf{v} = \mathbf{v}(\mathbf{u})$  serves as the Lagrange multiplier associated with the constraint  $\langle g \mathbf{m} \rangle = \mathbf{u}$  or equivalently under the closure assumption

$$(2.19) \quad \mathbf{u} = \langle \mathbf{m} \exp(\mathbf{v}(\mathbf{u}) \cdot \mathbf{m}) \rangle .$$

The moment system (2.16) can now be rewritten as

$$(2.20) \quad \mathbf{u}_{,t} + \mathbf{f}_{,x_i}^i = \mathbf{r}(\mathbf{u}) ,$$

where

$$(2.21) \quad \mathbf{f}^i = \langle v_i \mathbf{m} \exp(\mathbf{v}(\mathbf{u}) \cdot \mathbf{m}) \rangle .$$

Observe that using the kinetic Boltzmann structure, we have that

$$(2.22) \quad \mathcal{U}(\mathbf{v}) = \langle f \rangle = \langle \exp(\mathbf{v} \cdot \mathbf{m}) \rangle$$

is a suitable conjugate entropy function and that

$$(2.23) \quad U(\mathbf{u}) = \langle (\mathbf{v}(\mathbf{u}) \cdot \mathbf{m}) \exp(\mathbf{v}(\mathbf{u}) \cdot \mathbf{m}) - \exp(\mathbf{v}(\mathbf{u}) \cdot \mathbf{m}) \rangle$$

is the corresponding entropy function so that the duality relationship (2.9) holds.

The simplest example of a moment system is obtained by taking  $\mathbf{m}(v) = (1, v, |v|^2/2)^T$  corresponding to mass, momentum, and kinetic energy. In this instance, the collision integral vanishes identically,  $\mathbf{r} = 0$ , and (2.20) is the well-known system of Euler equations (5 moments) for a monotonic gas. More complex systems with 10, 14 or 35 moments have been considered in the literature [16]. In Appendix A, we give the corresponding Euler equations moment model for  $\gamma$ -law (polytropic) gases that is achieved by increasing the dimension of the phase space to include internal energy  $I$  and utilizing the moments  $\mathbf{m}(v, I) = (1, v, |v|^2/2 + I^\delta)^T$  for  $\delta = (1/(\gamma - 1) - d/2)^{-1}$ . In the case of the  $\gamma$ -law gas, one obtains a conjugate entropy function in  $\mathbb{R}^d$  of the form

$$(2.24) \quad \mathcal{U}(\mathbf{v}) = \langle c(\gamma, d) \exp(\mathbf{v} \cdot \mathbf{m}) \rangle , \quad c(\gamma, d) > 0$$

which is still compatible with the desired exponential structure. For brevity, we will omit constants such as  $c(\gamma, d)$  in our exponential form so that (2.22) may be regarded as an abstract form for (2.24) with suitably chosen phase space. From (2.22) it is clear that

$$(2.25) \quad \mathcal{U}_{, \mathbf{v}, \mathbf{v}} = \langle \mathbf{m} \otimes \mathbf{m} \exp(\mathbf{v} \cdot \mathbf{m}) \rangle$$

is SPD, i.e. the following double contraction to a scalar is positive

$$(2.26) \quad \mathcal{U}_{, \mathbf{v}_i, \mathbf{v}_j} \mathbf{z}_i \mathbf{z}_j = \langle (\mathbf{m} \cdot \mathbf{z})^2 \exp(\mathbf{v} \cdot \mathbf{m}) \rangle > 0 , \quad |\mathbf{z}| \neq 0 .$$

Furthermore,  $U_{, \mathbf{u}, \mathbf{u}} = \mathcal{U}_{, \mathbf{v}, \mathbf{v}}^{-1}$  is also SPD, hence  $U$  is also a convex function of  $\mathbf{u}$ . Consequently every system with the considered structure is hyperbolic symmetrizable and has a convex entropy  $U$  which is locally dissipated. This technique provides one of the simplest proofs of convexity for entropy functions associated with first order nonlinear conservation law systems derivable as moment closures of kinetic Boltzmann-like equations. In the case of the Euler equations of gasdynamics, the reader should compare this technique with the somewhat tedious proofs of convexity given in Refs. [12], [13], [9]. Finally, we mention the following general result for kinetic Boltzmann moment hierarchical systems which is used in later development.

**LEMMA 2.1. Generalized Convexity of Boltzmann Moment Conjugate Entropies.** *Let  $\mathbb{N} = \{0, 1, 2, \dots\}$  denote the set of nonnegative integers. All  $2k$  derivatives of the kinetic Boltzmann moment conjugate entropy (2.22)*

$$\mathcal{U}(\mathbf{v}) = \langle \exp(\mathbf{v} \cdot \mathbf{m}) \rangle$$

are SPD for  $k \in \mathbb{N}$

$$(2.27) \quad \frac{\partial^{2k} \mathcal{U}}{\partial \mathbf{v}^{2k}} \underbrace{[\mathbf{z}, \mathbf{z}, \dots, \mathbf{z}]}_{2k \text{ times}} > 0, \quad |\mathbf{z}| \neq 0.$$

**Proof.** Successive differentiation of (2.22)  $2k$  times yields the symmetric rank- $2k \times m$  tensor

$$(2.28) \quad \frac{\partial^{2k} \mathcal{U}}{\partial \mathbf{v}^{2k}}(\mathbf{v}) = \underbrace{(\mathbf{m} \otimes \mathbf{m} \otimes \dots \otimes \mathbf{m})}_{2k \text{ times}} \exp(\mathbf{v} \cdot \mathbf{m}),$$

followed by contraction to a scalar by a nonzero vector  $\mathbf{z} \in \mathbb{R}^m$

$$(2.29) \quad \frac{\partial^{2k} \mathcal{U}}{\partial \mathbf{v}^{2k}} \underbrace{[\mathbf{z}, \mathbf{z}, \dots, \mathbf{z}]}_{2k \text{ times}} = \langle (\mathbf{m} \cdot \mathbf{z})^{2k} \exp(\mathbf{v} \cdot \mathbf{m}) \rangle.$$

The moment vector  $\mathbf{m}$  contains  $m$  linearly independent polynomials spanning  $\mathbb{R}^m$ . The condition  $\mathbf{m}(v) \cdot \mathbf{z} = 0$  for fixed nonzero  $\mathbf{z}$  and variable  $v$  would violate the assumption of linear independence, thus we conclude that  $\mathbf{m}(v) \cdot \mathbf{z} \neq 0$  a.e., namely, except at points of measure zero in the phase space Lebesgue integration. The term  $\exp(\mathbf{v} \cdot \mathbf{m})$  is also positive for finite argument values in the phase space integration, hence we conclude for nonnegative powers  $2k$

$$(2.30) \quad \langle (\mathbf{m} \cdot \mathbf{z})^{2k} \exp(\mathbf{v} \cdot \mathbf{m}) \rangle > 0$$

and the stated lemma. ■

**2.2. The Eigenvector Scaling Theorem and Generalized Matrix Functions with Respect to the  $\tilde{A}_0$  Inner Product.** Next, we consider an important algebraic property of right symmetrizable systems which is used later in the implementation of the DG scheme. Simplifying upon the previous notation, let  $\tilde{A}_0 = \mathbf{u}, \mathbf{v}$ ,  $A_i = \mathbf{f}_{\mathbf{u}}^i$ ,  $\tilde{A}_i = A_i \tilde{A}_0$  and rewrite (2.3)

$$(2.31) \quad \tilde{A}_0 \mathbf{v}_{,t} + \tilde{A}_i \mathbf{v}_{,x_i} = 0.$$

The following theorem states a property of the symmetric matrix  $\tilde{A}_i$  symmetrized via the symmetric positive definite matrix  $\tilde{A}_0$ .

**THEOREM 2.2 (Eigenvector Scaling).** *Let  $A \in \mathbb{R}^{n \times n}$  be an arbitrary diagonalizable matrix and  $S$  the set of all right symmetrizers:*

$$S = \{B \in \mathbb{R}^{n \times n} \mid B \text{ SPD, } AB \text{ symmetric}\}.$$

*Further, let  $R \in \mathbb{R}^{n \times n}$  denote the right eigenvector matrix which diagonalizes  $A$*

$$A = R \Lambda R^{-1}$$

*with  $r$  distinct eigenvalues,  $\Lambda = \text{Diag}(\lambda_1 I_{m_1 \times m_1}, \lambda_2 I_{m_2 \times m_2}, \dots, \lambda_r I_{m_r \times m_r})$ . Then for each  $B \in S$  there exists a symmetric block diagonal matrix  $T = \text{Diag}(T_{m_1 \times m_1}, T_{m_2 \times m_2}, \dots, T_{m_r \times m_r})$  that block scales columns of  $R$ ,  $\tilde{R} = RT$ , such that*

$$B = \tilde{R} \tilde{R}^T, \quad A = \tilde{R} \Lambda \tilde{R}^{-1}$$

*which imply that*

$$AB = \tilde{R} \Lambda \tilde{R}^T.$$

**Proof.** Omitted, see [2]. ■

This last formula states a congruence relationship since  $\tilde{R}$  is not generally orthonormal and  $\Lambda$  does not represent the eigenvalues of  $AB$ . We shall refer to  $\tilde{R}$  as containing “entropy scaled” eigenvectors. Note that we can consider scalar combinations of  $\tilde{A}_i$  with the same scaling properties for arbitrary  $\mathbf{n} \in \mathbb{R}^m$ , i.e.

$$(2.32) \quad \tilde{A}(\mathbf{n}) \equiv \mathbf{n}_i \tilde{A}_i = \tilde{R}(\mathbf{n}) \Lambda(\mathbf{n}) \tilde{R}^T(\mathbf{n}) , \quad \tilde{A}_0 = \tilde{R}(\mathbf{n}) \tilde{R}^T(\mathbf{n}) .$$

Wavespeeds associated with the system (2.31) and the direction vector  $\mathbf{n}$  are given by critical values of the Rayleigh quotient

$$(2.33) \quad \frac{\xi^T \tilde{A}(\mathbf{n}) \xi}{\xi^T \tilde{A}_0 \xi} = \frac{\xi^T \tilde{R}(\mathbf{n}) \Lambda(\mathbf{n}) \tilde{R}^T(\mathbf{n}) \xi}{\xi \cdot \tilde{R}(\mathbf{n}) \tilde{R}^T(\mathbf{n}) \xi} = \frac{\eta^T \Lambda(\mathbf{n}) \eta}{\eta \cdot \eta} , \quad \xi, \eta \in \mathbb{R}^m, \quad \eta = \tilde{R} \xi, \quad |\xi| \neq 0 ,$$

which are simply elements of  $\Lambda(\mathbf{n})$ . For use in later developments, it is useful to define a matrix function  $f_{\tilde{A}_0}(\tilde{A})$  with respect to the Riemannian matrix  $\tilde{A}_0$  with critical values of the Rayleigh quotient given by  $f(\Lambda_{ii}), i = 1, \dots, m$ . This matrix function takes a particularly simple form as given by the following proposition:

**PROPOSITION 2.3.** *Barth [2, 1]. Let  $\tilde{A}_0$  denote the SPD right symmetrizer of  $A$  such that  $\tilde{A} = A\tilde{A}_0$ ,  $\tilde{A}_0 = \tilde{R}\tilde{R}^T$ , and  $A = \tilde{R}\Lambda\tilde{R}^{-1}$ . The generalized matrix function  $f_{\tilde{A}_0}(\tilde{A})$  is symmetric and defined canonically in terms of entropy scaled eigenvectors as*

$$(2.34) \quad f_{\tilde{A}_0}(\tilde{A}) = \tilde{R}f(\Lambda)\tilde{R}^T$$

where  $f(\Lambda)$  is performed componentwise.

**Proof.** Assume the desired critical values  $f(\Lambda)$  and the Rayleigh quotient producing them

$$(2.35) \quad \frac{\eta^T f(\Lambda) \eta}{\eta \cdot \eta} = \frac{\xi^T \tilde{R} f(\Lambda) \tilde{R}^T \xi}{\xi^T \tilde{R} \tilde{R}^T \xi} = \frac{\xi^T f_{\tilde{A}_0}(\tilde{A}) \xi}{\xi^T \tilde{A}_0 \xi} , \quad \xi, \eta \in \mathbb{R}^m, \quad \eta = \tilde{R} \xi, \quad |\xi| \neq 0 .$$

see also [2, 1]. ■

In later sections, the generalized matrix absolute value function  $|\tilde{A}|_{\tilde{A}_0}$  will be required. Using Proposition (2.3) stated above

$$(2.36) \quad |\tilde{A}|_{\tilde{A}_0} = \tilde{R}|\Lambda|\tilde{R}^T .$$

Finally, observe that using these scaled eigenvectors,  $\tilde{R}$ , we have the following equivalent representations<sup>1</sup> of  $\tilde{A}$  and  $\tilde{A}_0$  that are used in later developments:

$$(2.37) \quad \tilde{A} = \sum_{i=1}^m \lambda_i \tilde{r}_i \otimes \tilde{r}_i , \quad \tilde{A}_0 = \sum_{i=1}^m \tilde{r}_i \otimes \tilde{r}_i ,$$

and

$$(2.38) \quad |\tilde{A}|_{\tilde{A}_0} = \sum_{i=1}^m |\lambda_i| \tilde{r}_i \otimes \tilde{r}_i$$

where  $\tilde{r}_i$  denotes the  $i$ -th column of  $\tilde{R}$ .

<sup>1</sup>These representations should not be confused with the spectral decomposition of a matrix by orthonormal transform.

**3. DG Finite Element Method.** Let  $\Omega$  denote a spatial domain composed of nonoverlapping elements  $T_i$ ,  $\Omega = \cup T_i$ ,  $T_i \cap T_j = \emptyset$ ,  $i \neq j$  and  $I^n = ]t^n, t^{n+1}[$  the  $n$ -th time interval. It is useful to also define the element set  $\mathcal{T} = \{T_1, T_2, \dots, T_{|\mathcal{T}|}\}$  and edge set  $\mathcal{E} = \{e_1, e_2, \dots, e_{|\mathcal{E}|}\}$ . To simplify the exposition, consider a single variational formulation with weakly enforced boundary conditions. In the DG formulations (see [15, 3] and references therein), functions are discontinuous in space and time, i.e.

$$\mathcal{V}^h = \left\{ \mathbf{v}^h \mid \mathbf{v}|_{T \times I^n} \in \left( \mathcal{P}_k(T \times I^n) \right)^m \right\} .$$

For ease of exposition, we consider a spatial domain  $\Omega$  which is either periodic in all space dimensions or nonperiodic with compactly supported initial data. Consider the first order Cauchy system

$$(3.1) \quad \begin{cases} \mathbf{u}_t + \mathbf{f}_{,x_i}^i = 0 & \text{in } \Omega \\ \mathbf{u}(x, 0) = \mathbf{u}_0(x) \end{cases}$$

with  $A(\mathbf{n}) = \mathbf{n}_i A_i$  and  $\tilde{A}(\mathbf{n}) = \mathbf{n}_i \tilde{A}_i$ . The DG scheme with weakly imposed boundary conditions in time is defined by the following stabilized variational formulation:

Find  $\mathbf{v}^h \in \mathcal{V}^h$  such that for all  $\mathbf{w}^h \in \mathcal{V}^h$

$$(3.2) \quad B(\mathbf{v}^h, \mathbf{w}^h)_{\text{GAL}} = 0$$

where

$$\begin{aligned} B(\mathbf{v}, \mathbf{w})_{\text{GAL}} = & \int_{I^n} \int_{\Omega} (-\mathbf{u}(\mathbf{v}) \cdot \mathbf{w}_{,t} - \mathbf{f}^i(\mathbf{v}) \cdot \mathbf{w}_{,x_i}) dx dt \\ & + \int_{\Omega} (\mathbf{w}(t_-^{n+1}) \cdot \mathbf{u}(\mathbf{v}(t_-^{n+1})) - \mathbf{w}(t_+^n) \cdot \mathbf{u}(\mathbf{v}(t_+^n))) dx \\ & + \int_{I^n} \sum_{e \in \mathcal{E}} \int_e (\mathbf{w}(x_-) - \mathbf{w}(x_+)) \cdot \mathbf{h}(\mathbf{v}(x_-), \mathbf{v}(x_+); \mathbf{n}) dx dt \end{aligned}$$

where  $\mathbf{h}$  denotes a numerical flux function. Throughout, we consider numerical fluxes of the form

$$(3.3) \quad \mathbf{h}(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) = \frac{1}{2} (\mathbf{f}(\mathbf{v}_-; \mathbf{n}) + \mathbf{f}(\mathbf{v}_+; \mathbf{n})) - \frac{1}{2} \mathbf{h}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) .$$

These fluxes are consistent with the true flux in the sense that  $\mathbf{f}(\mathbf{v}; \mathbf{n}) = \mathbf{h}(\mathbf{v}, \mathbf{v}; \mathbf{n})$ .

**3.1. DG Nonlinear Energy Analysis.** Before presenting the nonlinear energy result, we recall some supporting corollaries concerning entropy function/flux jump identities at space-time slab interfaces. Note that throughout this section, we utilize the state-space parameterization

$$\bar{\mathbf{v}}(\theta) \equiv \mathbf{v}(x_-) + \theta [\mathbf{v}]_{x_-}^{x_+}$$

(similarly across time slab interfaces) for use in state-space path integrations and the interface averaging operator

$$\langle\langle \mathbf{v} \rangle\rangle_{x_-}^{x_+} \equiv \frac{\mathbf{v}(x_-) + \mathbf{v}(x_+)}{2} .$$

**LEMMA 3.1. Interface Jump Identities.** *Barth [2, 1] Let  $Z(\mathbf{u}), \mathcal{Z}(\mathbf{v}) : \mathbb{R}^m \mapsto \mathbb{R}$  be twice differentiable functions of their argument satisfying the duality relationship*

$$(3.4) \quad Z(\mathbf{u}) + \mathcal{Z}(\mathbf{v}) = \mathcal{Z}_{,\mathbf{v}} \mathbf{v} .$$

*The following jump identities hold across interfaces*

$$(3.5) \quad [Z]_{x_-}^{x_+} - [\mathcal{Z}_{,\mathbf{v}}]_{x_-}^{x_+} \mathbf{v}(x_+) + \int_0^1 (1 - \theta) [\mathbf{v}]_{x_-}^{x_+} \cdot \mathcal{Z}_{,\mathbf{v},\mathbf{v}}(\bar{\mathbf{v}}(\theta)) [\mathbf{v}]_{x_-}^{x_+} d\theta = 0$$

$$(3.6) \quad [Z]_{x_-}^{x_+} - [\mathcal{Z}_{,\mathbf{v}}]_{x_-}^{x_+} \mathbf{v}(x_-) - \int_0^1 \theta [\mathbf{v}]_{x_-}^{x_+} \cdot \mathcal{Z}_{,\mathbf{v},\mathbf{v}}(\bar{\mathbf{v}}(\theta)) [\mathbf{v}]_{x_-}^{x_+} d\theta = 0 .$$

**Proof.** Omitted, see [2, 1]. ■

**COROLLARY 3.2. Temporal Space-Time Slab Interface Identity.** *Barth [2, 1]. Let  $t_{\pm}$  denote a temporal space-time slab interface. The following entropy function jump identity holds across time slab interfaces*

$$(3.7) \quad \int_{\Omega} \left( [U]_{t_{\pm}}^{t_{\pm}} - \mathbf{v}^T(t_{\pm}) [\mathbf{u}]_{t_{\pm}}^{t_{\pm}} \right) dx + \frac{1}{2} ||| [\mathbf{v}]_{t_{\pm}}^{t_{\pm}} |||_{\tilde{A}_0, \Omega}^2 = 0$$

where

$$(3.8) \quad ||| [\mathbf{v}]_{t_{\pm}}^{t_{\pm}} |||_{\tilde{A}_0, \Omega}^2 \equiv \int_{\Omega} \int_0^1 2(1-\theta) [\mathbf{v}]_{t_{\pm}}^{t_{\pm}} \cdot \tilde{A}_0(\bar{\mathbf{v}}(\theta)) [\mathbf{v}]_{t_{\pm}}^{t_{\pm}} d\theta dx \geq 0 .$$

**COROLLARY 3.3. Spatial Space-Time Slab Interface Identity.** *Barth [2, 1]. Let  $x_{\pm}$  denote a spatial element interface. The following entropy flux jump identity holds across spatial element interfaces*

$$(3.9) \quad [F^i]_{x_{\pm}}^{x_{\pm}} - \langle \langle \mathbf{v}^T \rangle \rangle_{x_{\pm}}^{x_{\pm}} [\mathbf{f}^i]_{x_{\pm}}^{x_{\pm}} + \frac{1}{2} \int_0^1 (1-2\theta) [\mathbf{v}]_{x_{\pm}}^{x_{\pm}} \cdot \tilde{A}_i(\bar{\mathbf{v}}(\theta)) [\mathbf{v}]_{x_{\pm}}^{x_{\pm}} d\theta = 0 .$$

Note that in actual numerical calculations, it is desirable to use the variational form given by (3.2) since integration by parts has been used to insure exact discrete conservation even with inexact numerical quadrature of the various integrals. For analysis purposes, however, it is desirable to use the following equivalent non-integrated-by-parts formulation:

Find  $\mathbf{v}^h \in \mathcal{V}^h$  such that for all  $\mathbf{w}^h \in \mathcal{V}^h$

$$(3.10) \quad B(\mathbf{v}^h, \mathbf{w}^h)_{\text{GAL}} = 0$$

where

$$\begin{aligned} B(\mathbf{v}, \mathbf{w})_{\text{GAL}} &= \int_{I^n} \int_{\Omega} \mathbf{w} \cdot (\mathbf{u}_{,t} + \mathbf{f}_{,x_i}^i(\mathbf{v})) dx dt \\ &\quad + \int_{\Omega} \mathbf{w}(t_+^n) \cdot [\mathbf{u}]_{t_+^n}^{t_+^n} dx \\ &\quad + \int_{I^n} \sum_{e \in \mathcal{E}} \int_e \frac{1}{2} [\mathbf{w}]_{x_{\pm}}^{x_{\pm}} \cdot \mathbf{h}^d(\mathbf{v}(x_-), \mathbf{v}(x_+); \mathbf{n}) dx dt \\ &\quad + \int_{I^n} \sum_{e \in \mathcal{E}} \int_e \langle \langle \mathbf{w} \rangle \rangle_{x_{\pm}}^{x_{\pm}} \cdot [\mathbf{f}(\mathbf{v}; \mathbf{n})]_{x_{\pm}}^{x_{\pm}} dx dt \end{aligned}$$

where  $\mathbf{h}^d$  denotes the flux dissipation term incorporated into the total numerical flux.

**THEOREM 3.4. DG Global Entropy Norm Stability (Nonlinear Hyperbolic System).** *The variational formulation (3.10) for nonlinear systems of conservation laws with convex entropy extension and symmetric mean-value flux dissipation*

$$\mathbf{h}_{\text{SMV}}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) = |\tilde{A}|_{\text{SMV}} [\mathbf{v}]_{x_{\pm}}^{x_{\pm}} , \quad |\tilde{A}|_{\text{SMV}} \equiv \int_0^1 |\tilde{A}(\bar{\mathbf{v}}(\theta); \mathbf{n})|_{\tilde{A}_0} d\theta$$

is entropy norm stable with the following global balance:

$$(3.11) \quad \frac{1}{2} \sum_{n=0}^{N-1} \left( ||| [\mathbf{v}]_{t_{\pm}}^{t_{\pm}} |||_{\tilde{A}_0, \Omega}^2 + \sum_{e \in \mathcal{E}} \langle [\mathbf{v}]_{x_{\pm}}^{x_{\pm}} \rangle_{|\tilde{A}|, e \times I^n}^2 \right) + \int_{\Omega} U(t_-^N) dx = \int_{\Omega} U(t_-^0) dx$$

with

$$|\tilde{A}(\mathbf{n})| = \int_0^1 2(1-\theta) \left( \tilde{A}^+(\bar{\mathbf{v}}(\theta); \mathbf{n})_{\tilde{A}_0} - \tilde{A}^-(\bar{\mathbf{v}}(1-\theta); \mathbf{n})_{\tilde{A}_0} \right) d\theta .$$



**Proof.** Repeated here from [2, 1]. Construct the energy balance for the interval  $[t_-^0, t_-^N] = \cup_{n=0}^{N-1} I^n$  by setting  $\mathbf{w} = \mathbf{v}$  and evaluating the various integrals. Consider the time derivative integral

$$\int_{\Omega} \int_{I^n} \mathbf{v}^T \mathbf{u}_{,t} dt dx = \int_{\Omega} \int_{I^n} U_{,t} dt dx = \int_{\Omega} \left( [U]_{t_-^n}^{t_-^{n+1}} - [U]_{t_-^n}^{t_-^n} \right) dx$$

and combine with the jump integral across time slabs. From Corollary 3.2

$$\int_{\Omega} \int_{I^n} \mathbf{v}^T \mathbf{u}_{,t} dt dx + \int_{\Omega} \mathbf{v}^T (t_+^n) [\mathbf{u}]_{t_-^n}^{t_+^n} dx = \int_{\Omega} [U]_{t_-^n}^{t_-^{n+1}} dx + \frac{1}{2} ||| [\mathbf{v}]_{t_-^n}^{t_+^n} |||_{\tilde{A}_0, \Omega}^2 .$$

When summed over all time slabs, the first term on the right-hand-side of this equation vanishes except for initial and final time slab contributions. Next, consider the spatial operator term and apply the divergence theorem

$$(3.12) \quad \int_{I^n} \int_{\Omega} \mathbf{v}^T \mathbf{f}_{,x_i} dx dt = \int_{I^n} \int_{\Omega} F_{,x_i}^i dx dt = \int_{I^n} \sum_{e \in \mathcal{E}} \int_e -[F(\mathbf{v}; \mathbf{n})]_{x_-^e}^{x_+^e} dx dt$$

where  $F(\mathbf{v}; \mathbf{n}) = \mathbf{n}_i F^i(\mathbf{v})$ . Combining all the space terms and applying Corollary 3.3

$$\begin{aligned} II_{\text{space}}^n &\equiv \int_{I^n} \sum_{e \in \mathcal{E}} \int_e \left( -[F(\mathbf{v}; \mathbf{n})]_{x_-^e}^{x_+^e} + \langle \mathbf{v} \rangle_{x_-^e}^{x_+^e} \cdot [\mathbf{f}(\mathbf{n})]_{x_-^e}^{x_+^e} + \frac{1}{2} [\mathbf{v}]_{x_-^e}^{x_+^e} \cdot \mathbf{h}^d \right) dx dt \\ &= \int_{I^n} \sum_{e \in \mathcal{E}} \int_e \frac{1}{2} [\mathbf{v}]_{x_-^e}^{x_+^e} \cdot \left( \mathbf{h}^d + \int_0^1 (1 - 2\theta) \tilde{A}_i(\bar{\mathbf{v}}(\theta)) [\mathbf{v}]_{x_-^e}^{x_+^e} d\theta \right) dx dt . \end{aligned}$$

In summary, collecting all terms and summing over time slabs we have

$$\begin{aligned} B(\mathbf{v}, \mathbf{v})_{\text{GAL}} &= \sum_{n=0}^{N-1} \left( \int_{\Omega} [U]_{t_-^n}^{t_-^{n+1}} dx + \frac{1}{2} ||| [\mathbf{v}]_{t_-^n}^{t_+^n} |||_{\tilde{A}_0, \Omega}^2 + II_{\text{space}}^n \right) \\ &= \int_{\Omega} U(t_-^N) dx - \int_{\Omega} U(t_-^0) dx + \sum_{n=0}^{N-1} \left( \frac{1}{2} ||| [\mathbf{v}]_{t_-^n}^{t_+^n} |||_{\tilde{A}_0, \Omega}^2 + II_{\text{space}}^n \right) . \end{aligned}$$

When written in this form, it becomes clear that a sufficient condition for energy stability is that for all time intervals  $I^n$

$$(3.13) \quad II_{\text{space}}^n \geq 0$$

which serves as a design condition for the flux dissipation.

$$\begin{aligned} II_{\text{space}}^n &= \int_{I^n} \sum_{e \in \mathcal{E}} \int_e \frac{1}{2} [\mathbf{v}]_{x_-^e}^{x_+^e} \cdot \left( \mathbf{h}^d + \int_0^1 (1 - 2\theta) \tilde{A}_i(\bar{\mathbf{v}}(\theta)) [\mathbf{v}]_{x_-^e}^{x_+^e} d\theta \right) dx dt \\ &= \int_{I^n} \sum_{e \in \mathcal{E}} \int_e \frac{1}{2} [\mathbf{v}]_{x_-^e}^{x_+^e} \cdot \left( \mathbf{h}^d + \int_0^1 (1 - \theta) \tilde{A}_i(\bar{\mathbf{v}}(\theta)) [\mathbf{v}]_{x_-^e}^{x_+^e} d\theta \right) dx dt \\ &\quad - \int_{I^n} \sum_{e \in \mathcal{E}} \int_e \frac{1}{2} [\mathbf{v}]_{x_-^e}^{x_+^e} \cdot \int_0^1 \theta \tilde{A}_i(\bar{\mathbf{v}}(\theta)) [\mathbf{v}]_{x_-^e}^{x_+^e} d\theta dx dt \end{aligned}$$

The choice

$$(3.14) \quad \mathbf{h}^d = \mathbf{h}_{\text{SMV}}^d \equiv \int_0^1 |\tilde{A}(\bar{\mathbf{v}}(\theta); \mathbf{n})|_{\tilde{A}_0} [\mathbf{v}]_{x_-^e}^{x_+^e} d\theta$$

yields

$$II_{\text{space}}^n = \int_{I^n} \sum_{e \in \mathcal{E}} \int_e \frac{1}{2} [\mathbf{v}]_{x_-^e}^{x_+^e} \cdot 2 \int_0^1 \left( (1 - \theta) \tilde{A}^+(\bar{\mathbf{v}}(\theta); \mathbf{n})_{\tilde{A}_0} - \theta \tilde{A}^-(\bar{\mathbf{v}}(\theta); \mathbf{n})_{\tilde{A}_0} \right) d\theta [\mathbf{v}]_{x_-^e}^{x_+^e} dx dt$$

$$\begin{aligned}
&= \int_{I^n} \sum_{e \in \mathcal{E}} \int_e \frac{1}{2} [\mathbf{v}]_{x_-}^{x_+} \cdot 2 \int_0^1 (1 - \theta) \left( \tilde{A}^+(\bar{\mathbf{v}}(\theta); \mathbf{n})_{\tilde{A}_0} - \tilde{A}^-(\bar{\mathbf{v}}(1 - \theta); \mathbf{n})_{\tilde{A}_0} \right) d\theta [\mathbf{v}]_{x_-}^{x_+} dx dt \\
&= \int_{I^n} \sum_{e \in \mathcal{E}} \int_e \frac{1}{2} [\mathbf{v}]_{x_-}^{x_+} \cdot |\tilde{A}(\mathbf{n})| [\mathbf{v}]_{x_-}^{x_+} dx dt \geq 0 .
\end{aligned}$$

This completes the sufficient condition for energy decay in time. ■

An important observation to be made concerning (3.13) and energy stability is that any flux dissipation  $\mathbf{h}^d$  for which

$$(3.15) \quad [\mathbf{v}]_{-}^{\pm} \cdot \mathbf{h}_{\text{SMV}}^d \leq [\mathbf{v}]_{-}^{\pm} \cdot \mathbf{h}^d$$

is also energy stable with increased energy decay. This observation is used in later sections in designing simplified flux functions for the DG method.

### 3.2. DG Nonlinear Energy Analysis for Kinetic Boltzmann Moment Closure Hierarchies.

In this section, we analyze nonlinear energy properties of the discontinuous Galerkin method assuming the kinetic Boltzmann moment closure structure discussed in Sect. 2.1. Recall that in this framework

$$(3.16) \quad \mathbf{u} = \langle \mathbf{m} \exp(\mathbf{v} \cdot \mathbf{m}) \rangle, \quad \mathbf{f}^i = \langle v; \mathbf{m} \exp(\mathbf{v} \cdot \mathbf{m}) \rangle$$

with conjugate entropy given by

$$(3.17) \quad \mathcal{U}(\mathbf{v}) = \langle \exp(\mathbf{v} \cdot \mathbf{m}) \rangle .$$

Using these definitions, we briefly examine energy stability of the DG method for these moment closure hierarchies.

**THEOREM 3.5. DG Global Entropy Norm Stability of Kinetic Boltzmann Moment Closure Hierarchies.** *The variational formulation (3.10) for nonlinear systems of conservation laws with convex entropy extension and kinetic Boltzmann moment closure symmetric mean-value flux dissipation*

$$\mathbf{h}_{\text{KSMV}}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) = |\tilde{A}|_{\text{KSMV}} [\mathbf{v}]_{x_-}^{x_+}, \quad |\tilde{A}|_{\text{KSMV}} \equiv \int_0^1 \langle |v \cdot \mathbf{n}| \mathbf{m} \otimes \mathbf{m} \exp(\bar{\mathbf{v}}(\theta) \cdot \mathbf{m}) \rangle d\theta$$

is entropy norm stable with the following global balance:

$$(3.18) \quad \frac{1}{2} \sum_{n=0}^{N-1} \left( ||| [\mathbf{v}]_{t_-}^{t_+} |||_{\tilde{A}_0, \Omega}^2 + \sum_{e \in \mathcal{E}} \langle ([\mathbf{v}]_{x_-}^{x_+})^2_{|\tilde{A}|, e \times I^n} \rangle \right) + \int_{\Omega} U(t^N) dx = \int_{\Omega} U(t^0) dx$$

with

$$|\tilde{A}(\mathbf{n})| = \int_0^1 2(1 - \theta) \left( \tilde{A}^+(\bar{\mathbf{v}}(\theta); \mathbf{n})_{\tilde{A}_0} - \tilde{A}^-(\bar{\mathbf{v}}(1 - \theta); \mathbf{n})_{\tilde{A}_0} \right) d\theta .$$

**Proof.** By making the following generalizations

$$U(\mathbf{u}) = \langle (\mathbf{v}(\mathbf{u}) \cdot \mathbf{m}) \exp(\mathbf{v}(\mathbf{u}) \cdot \mathbf{m}) - \exp(\mathbf{v}(\mathbf{u}) \cdot \mathbf{m}) \rangle$$

$$\tilde{A}_0(\mathbf{v}) = \langle \mathbf{m} \otimes \mathbf{m} \exp(\mathbf{v} \cdot \mathbf{m}) \rangle$$

$$\tilde{A}(\mathbf{v}; \mathbf{n}) = \langle (v \cdot \mathbf{n}) \mathbf{m} \otimes \mathbf{m} \exp(\mathbf{v} \cdot \mathbf{m}) \rangle$$

$$\tilde{A}^{\pm}(\mathbf{v}; \mathbf{n})_{\tilde{A}_0} = \langle (v \cdot \mathbf{n})^{\pm} \mathbf{m} \otimes \mathbf{m} \exp(\mathbf{v} \cdot \mathbf{m}) \rangle ,$$

we can appeal once again to the space-time slab jump identities stated in lemma 3.1 and corollaries 3.7 and 3.3. Using these definitions and results, the proof of theorem 3.4 applies without alteration up to and including the equation

$$\begin{aligned} II_{\text{space}}^n &= \int_{I^n} \sum_{e \in \mathcal{E}} \int_e \frac{1}{2} [\mathbf{v}]_{x-}^{x+} \cdot \left( \mathbf{h}^d + \int_0^1 (1-2\theta) \tilde{A}_i(\bar{\mathbf{v}}(\theta)) [\mathbf{v}]_{x-}^{x+} d\theta \right) dx dt \\ &= \int_{I^n} \sum_{e \in \mathcal{E}} \int_e \frac{1}{2} [\mathbf{v}]_{x-}^{x+} \cdot \left( \mathbf{h}^d + \int_0^1 (1-\theta) \tilde{A}_i(\bar{\mathbf{v}}(\theta)) [\mathbf{v}]_{x-}^{x+} d\theta \right) dx dt \\ &\quad - \int_{I^n} \sum_{e \in \mathcal{E}} \int_e \frac{1}{2} [\mathbf{v}]_{x-}^{x+} \cdot \int_0^1 \theta \tilde{A}_i(\bar{\mathbf{v}}(\theta)) [\mathbf{v}]_{x-}^{x+} d\theta dx dt \end{aligned}$$

and the design condition

$$(3.19) \quad II_{\text{space}}^n \geq 0 .$$

In the present case, the choice

$$(3.20) \quad \mathbf{h}^d = \mathbf{h}_{\text{KSMV}}^d \equiv \int_0^1 \langle |v \cdot \mathbf{n}| \mathbf{m} \otimes \mathbf{m} \exp(\bar{\mathbf{v}}(\theta) \cdot \mathbf{m}) \rangle d\theta [\mathbf{v}]_{x-}^{x+}$$

is sufficient

$$\begin{aligned} II_{\text{space}}^n &= \int_{I^n} \sum_{e \in \mathcal{E}} \int_e \frac{1}{2} [\mathbf{v}]_{x-}^{x+} \cdot 2 \int_0^1 \left( (1-\theta) \tilde{A}^+(\bar{\mathbf{v}}(\theta); \mathbf{n})_{\tilde{A}_0} - \theta \tilde{A}^-(\bar{\mathbf{v}}(\theta); \mathbf{n})_{\tilde{A}_0} \right) d\theta [\mathbf{v}]_{x-}^{x+} dx dt \\ &= \int_{I^n} \sum_{e \in \mathcal{E}} \int_e \frac{1}{2} [\mathbf{v}]_{x-}^{x+} \cdot 2 \int_0^1 (1-\theta) \left( \tilde{A}^+(\bar{\mathbf{v}}(\theta); \mathbf{n})_{\tilde{A}_0} - \tilde{A}^-(\bar{\mathbf{v}}(1-\theta); \mathbf{n})_{\tilde{A}_0} \right) d\theta [\mathbf{v}]_{x-}^{x+} dx dt \\ &= \int_{I^n} \sum_{e \in \mathcal{E}} \int_e \frac{1}{2} [\mathbf{v}]_{x-}^{x+} \cdot |\tilde{A}(\mathbf{n})| [\mathbf{v}]_{x-}^{x+} dx dt \geq 0 . \end{aligned}$$

This completes the sufficient condition for energy decay in time for kinetic Boltzmann moment closure hierarchies. ■

**4. Simplified Numerical Flux Formulas for the DG Method.** The theoretical results of Sect. 3 provide the framework for constructing, analyzing, and proving energy stability for a number of simplified numerical flux functions. This task is undertaken in the remainder of this section. We are unaware of any previous DG analysis for systems ( $m > 1$ ) of nonlinear conservation laws which rigorously establishes energy stability for the fluxes considered here. Throughout, we use the notation  $\tilde{A} = \tilde{R} \Lambda \tilde{R}^T$  as defined earlier with  $\Lambda \equiv \text{diag}(\lambda_1, \dots, \lambda_m)$  assuming ordered entries  $\lambda_1 \leq \dots \leq \lambda_m$ . For numerical fluxes such as the SHHLE and SHLLEM flux, we also require that  $\lambda_1$  and  $\lambda_m$  be distinct in order that the construction be well defined.

• **Symmetric Lax-Friedrichs Flux (SLF)**

$$(4.1) \quad \mathbf{h}_{\text{SLF}}(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) = \frac{1}{2} (\mathbf{f}(\mathbf{v}_-; \mathbf{n}) + \mathbf{f}(\mathbf{v}_+; \mathbf{n})) - \frac{1}{2} \lambda_{\max} [\mathbf{u}(\mathbf{v})]_{x-}^{x+}$$

with

$$\lambda_{\max} \equiv \sup_{0 \leq \theta \leq 1} (\lambda_m(\bar{\mathbf{v}}(\theta))) .$$

• **Symmetric Lax-Friedrichs Matrix Flux (SLFM)**

$$(4.2) \quad \mathbf{h}_{\text{SLFM}}(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) = \frac{1}{2} (\mathbf{f}(\mathbf{v}_-; \mathbf{n}) + \mathbf{f}(\mathbf{v}_+; \mathbf{n})) - \frac{1}{2} \mathbf{h}_{\text{SLFM}}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n})$$

with

$$\mathbf{h}_{\text{SLFM}}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) = \lambda_{\max} \frac{1}{2} (\tilde{A}_0(\mathbf{v}_-) + \tilde{A}_0(\mathbf{v}_+)) [\mathbf{v}]_{x-}^{x+}$$

and

$$\lambda_{\max} \equiv \sup_{0 \leq \theta \leq 1} (\lambda_m(\bar{v}(\theta)))$$

- **Symmetric Harten-Lax-van Leer Flux (SHLLE)**

$$(4.3) \quad \mathbf{h}_{\text{SHLLE}}(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) = \frac{1}{2} (\mathbf{f}(\mathbf{v}_-; \mathbf{n}) + \mathbf{f}(\mathbf{v}_+; \mathbf{n})) - \frac{1}{2} \mathbf{h}_{\text{SHLLE}}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n})$$

with

$$\mathbf{h}_{\text{SHLLE}}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) = \frac{\lambda_{\max} + \lambda_{\min}}{\lambda_{\max} - \lambda_{\min}} [\mathbf{f}(\mathbf{v}; \mathbf{n})]_{-}^{+} - \frac{2\lambda_{\max}\lambda_{\min}}{\lambda_{\max} - \lambda_{\min}} [\mathbf{u}(\mathbf{v})]_{-}^{+} .$$

and

$$\lambda_{\max} = \sup_{0 \leq \theta \leq 1} \max(0, \lambda_m(\bar{v}(\theta))) , \quad \lambda_{\min} = \inf_{0 \leq \theta \leq 1} \min(0, \lambda_1(\bar{v}(\theta))) .$$

- **Symmetric Harten-Lax-van Leer Modified Flux (SHLLEM)**

$$(4.4) \quad \mathbf{h}_{\text{SHLLEM}}(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) = \frac{1}{2} (\mathbf{f}(\mathbf{v}_-; \mathbf{n}) + \mathbf{f}(\mathbf{v}_+; \mathbf{n})) - \frac{1}{2} \mathbf{h}_{\text{SHLLEM}}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n})$$

with

$$\mathbf{h}_{\text{SHLLEM}}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) = \frac{\lambda_{\max} + \lambda_{\min}}{\lambda_{\max} - \lambda_{\min}} [\mathbf{f}(\mathbf{v}; \mathbf{n})]_{-}^{+} - \frac{2\lambda_{\max}\lambda_{\min}}{\lambda_{\max} - \lambda_{\min}} [\mathbf{u}(\mathbf{v})]_{-}^{+} - \sum_{i=2}^{m-1} \int_0^1 f_i(\theta) \tilde{\tau}_i \otimes \tilde{\tau}_i [\mathbf{v}]_{-}^{+} d\theta$$

with

$$f_i(\theta) = \frac{\max(0, \lambda_m(\theta)) + \min(0, \lambda_1(\theta))}{\max(0, \lambda_m(\theta)) - \min(0, \lambda_1(\theta))} \lambda_i(\theta) - \frac{2 \max(0, \lambda_m(\theta)) \min(0, \lambda_1(\theta))}{\max(0, \lambda_m(\theta)) - \min(0, \lambda_1(\theta))} - |\lambda_i(\theta)|$$

and

$$\lambda_{\max} = \sup_{0 \leq \theta \leq 1} \max(0, \lambda_m(\bar{v}(\theta))) , \quad \lambda_{\min} = \inf_{0 \leq \theta \leq 1} \min(0, \lambda_1(\bar{v}(\theta))) .$$

- **Discrete Kinetic Symmetric Mean-Value Flux (DKSMV)**

$$(4.5) \quad \mathbf{h}_{\text{DKSMV}}(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) = \frac{1}{2} (\mathbf{f}(\mathbf{v}_-; \mathbf{n}) + \mathbf{f}(\mathbf{v}_+; \mathbf{n})) - \frac{1}{2} \mathbf{h}_{\text{DKSMV}}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n})$$

with

$$\begin{aligned} \mathbf{h}_{\text{DKSMV}}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) &= \frac{1}{2} \langle |v \cdot \mathbf{n}| \mathbf{m} \otimes \mathbf{m} (\exp(\mathbf{v}_- \cdot \mathbf{m}(v)) + \exp(\mathbf{v}_+ \cdot \mathbf{m}(v))) \rangle [\mathbf{v}]_{-}^{+} \\ &= \frac{1}{2} \langle |v \cdot \mathbf{n}| \mathbf{m} \otimes \mathbf{m} \exp(\mathbf{v}_- \cdot \mathbf{m}(v)) \rangle [\mathbf{v}]_{-}^{+} \\ &\quad + \frac{1}{2} \langle |v \cdot \mathbf{n}| \mathbf{m} \otimes \mathbf{m} \exp(\mathbf{v}_+ \cdot \mathbf{m}(v)) \rangle [\mathbf{v}]_{-}^{+} \end{aligned}$$

Observe that explicit path  $\theta$ -integration has been avoided in all these simplified fluxes (except correction terms in SHLLEM). In addition, we have the following theorem:

**THEOREM 4.1. Energy Stability of Simplified Flux Formulas** *The variational formulation (3.10) for nonlinear systems of conservation laws utilizing any of the candidate approximate numerical fluxes (4.1), (4.2), (4.3), (4.4), (4.5) is entropy norm stable in the sense of Theorem 3.4 or 3.5.*

**Proofs:** Given on a cases-by-case basis.

Symmetric Lax-Friedrichs Flux (SLF):

$$\begin{aligned}
[\mathbf{v}]_{x-}^{x+} \cdot \mathbf{h}_{\text{SMV}}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) &= [\mathbf{v}]_{x-}^{x+} \cdot \int_0^1 |\tilde{A}(\bar{\mathbf{v}}(\theta); \mathbf{n})|_{\tilde{A}_0} [\mathbf{v}]_{x-}^{x+} d\theta \\
&= \int_0^1 [\mathbf{v}]_{x-}^{x+} \cdot \tilde{R}(\bar{\mathbf{v}}(\theta); \mathbf{n}) |\Lambda(\bar{\mathbf{v}}(\theta))| \tilde{R}^T(\bar{\mathbf{v}}(\theta); \mathbf{n}) [\mathbf{v}]_{x-}^{x+} d\theta \\
&\leq \sup_{0 \leq \theta \leq 1} (\lambda_m(\bar{\mathbf{v}}(\theta))) \int_0^1 [\mathbf{v}]_{x-}^{x+} \cdot \tilde{R}(\bar{\mathbf{v}}(\theta); \mathbf{n}) \tilde{R}^T(\bar{\mathbf{v}}(\theta); \mathbf{n}) [\mathbf{v}]_{x-}^{x+} d\theta \\
&= \sup_{0 \leq \theta \leq 1} (\lambda_m(\bar{\mathbf{v}}(\theta))) \int_0^1 [\mathbf{v}]_{x-}^{x+} \cdot \tilde{A}_0(\bar{\mathbf{v}}(\theta)) [\mathbf{v}]_{x-}^{x+} d\theta \\
&= \sup_{0 \leq \theta \leq 1} (\lambda_m(\bar{\mathbf{v}}(\theta))) [\mathbf{v}]_{x-}^{x+} \cdot [\mathbf{u}(\mathbf{v})]_{x-}^{x+} \\
&= [\mathbf{v}]_{x-}^{x+} \cdot \mathbf{h}_{\text{SLF}}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) \quad \blacksquare
\end{aligned}$$

Symmetric Lax-Friedrichs Matrix Flux (SLFM):

$$\begin{aligned}
[\mathbf{v}]_{x-}^{x+} \cdot \mathbf{h}_{\text{SMV}}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) &= [\mathbf{v}]_{x-}^{x+} \cdot \int_0^1 |\tilde{A}(\bar{\mathbf{v}}(\theta); \mathbf{n})|_{\tilde{A}_0} [\mathbf{v}]_{x-}^{x+} d\theta \\
&= \int_0^1 [\mathbf{v}]_{x-}^{x+} \cdot \tilde{R}(\bar{\mathbf{v}}(\theta); \mathbf{n}) |\Lambda(\bar{\mathbf{v}}(\theta))| \tilde{R}^T(\bar{\mathbf{v}}(\theta); \mathbf{n}) [\mathbf{v}]_{x-}^{x+} d\theta \\
&\leq \sup_{0 \leq \theta \leq 1} (\lambda_m(\bar{\mathbf{v}}(\theta))) \int_0^1 [\mathbf{v}]_{x-}^{x+} \cdot \tilde{R}(\bar{\mathbf{v}}(\theta); \mathbf{n}) \tilde{R}^T(\bar{\mathbf{v}}(\theta); \mathbf{n}) [\mathbf{v}]_{x-}^{x+} d\theta \\
&= \sup_{0 \leq \theta \leq 1} (\lambda_m(\bar{\mathbf{v}}(\theta))) \int_0^1 [\mathbf{v}]_{x-}^{x+} \cdot \tilde{A}_0(\bar{\mathbf{v}}(\theta)) [\mathbf{v}]_{x-}^{x+} d\theta
\end{aligned}$$

Examining the scalar function

$$g(\theta) \equiv [\mathbf{v}]_{x-}^{x+} \cdot \tilde{A}_0(\bar{\mathbf{v}}(\theta)) [\mathbf{v}]_{x-}^{x+} = [\mathbf{v}]_{x-}^{x+} \cdot \frac{\partial^2 \mathcal{U}}{\partial \mathbf{v}^2}(\bar{\mathbf{v}}(\theta)) [\mathbf{v}]_{x-}^{x+} ,$$

differentiation yields

$$g'' = \frac{\partial^4 \mathcal{U}}{\partial \mathbf{v}^4}(\bar{\mathbf{v}}(\theta)) [[\mathbf{v}]_{x-}^{x+}, [\mathbf{v}]_{x-}^{x+}, [\mathbf{v}]_{x-}^{x+}, [\mathbf{v}]_{x-}^{x+}]$$

where the right-hand-side term denotes the rank-4 contraction to a scalar. But for systems derived as moments of the kinetic Boltzmann equation, we have from Sect. 2.1 that

$$\frac{\partial^4 \mathcal{U}}{\partial \mathbf{v}^4}[\mathbf{z}, \mathbf{z}, \mathbf{z}, \mathbf{z}] > 0, \quad |\mathbf{z}| \neq 0 .$$

Consequently,  $g(\theta)$  is convex for all  $\theta$ ,

$$g(\theta) \leq (1 - \theta) g(0) + \theta g(1) ,$$

so that

$$\begin{aligned}
[\mathbf{v}]_{x-}^{x+} \cdot \int_0^1 [\mathbf{v}]_{x-}^{x+} \cdot \tilde{A}_0(\bar{\mathbf{v}}(\theta)) [\mathbf{v}]_{x-}^{x+} d\theta &\leq \frac{1}{2} \left( [\mathbf{v}]_{x-}^{x+} \cdot \tilde{A}_0(\mathbf{v}_-) [\mathbf{v}]_{x-}^{x+} + [\mathbf{v}]_{x-}^{x+} \cdot \tilde{A}_0(\mathbf{v}_+) [\mathbf{v}]_{x-}^{x+} \right) \\
&= \frac{1}{2} [\mathbf{v}]_{x-}^{x+} \cdot \left( \tilde{A}_0(\mathbf{v}_-) + \tilde{A}_0(\mathbf{v}_+) \right) [\mathbf{v}]_{x-}^{x+} .
\end{aligned}$$

This yields

$$[\mathbf{v}]_{x-}^{x+} \cdot \mathbf{h}_{\text{SMV}}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) \leq \sup_{0 \leq \theta \leq 1} (\lambda_m(\bar{\mathbf{v}}(\theta))) \frac{1}{2} [\mathbf{v}]_{x-}^{x+} \cdot \left( \tilde{A}_0(\mathbf{v}_-) + \tilde{A}_0(\mathbf{v}_+) \right) [\mathbf{v}]_{x-}^{x+}$$

$$= [\mathbf{v}]_{x_-}^{z+} \cdot \mathbf{h}_{\text{SLFM}}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) \quad \blacksquare$$

Symmetric Harten-Lax-van Leer Fluxes (SHLLE) and (SHLLEM):

Our first task will be to prove that

$$[\mathbf{v}]_-^+ \cdot \mathbf{h}_{\text{SMV}}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) \leq [\mathbf{v}]_-^+ \cdot \mathbf{h}_{\text{SHLLEM}}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n})$$

followed by

$$[\mathbf{v}]_-^+ \cdot \mathbf{h}_{\text{SHLLEM}}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) \leq [\mathbf{v}]_-^+ \cdot \mathbf{h}_{\text{SHLLE}}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) .$$

To do so, consider the symmetrization matrices  $\tilde{A}$  and  $\tilde{A}_0$  written in the form (2.37)

$$\tilde{A} = \sum_{i=1}^m \lambda_i \tilde{r}_i \otimes \tilde{r}_i \quad \tilde{A}_0 = \sum_{i=1}^m \tilde{r}_i \otimes \tilde{r}_i$$

and the following useful  $2 \times 2$  matrix identity for  $\tilde{A} \in \mathbb{R}^{2 \times 2}$  and  $\tilde{A}_0 \in \mathbb{R}^{2 \times 2}$

$$(4.6) \quad |\tilde{A}|_{\tilde{A}_0} = \frac{\max(0, \lambda_2) + \min(0, \lambda_1)}{\max(0, \lambda_2) - \min(0, \lambda_1)} \tilde{A} - \frac{2 \max(0, \lambda_2) \min(0, \lambda_1)}{\max(0, \lambda_2) - \min(0, \lambda_1)} \tilde{A}_0$$

as can be easily verified by substitution. Generalizing to  $m \times m$  matrices with  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_m$  and  $\lambda_1 < \lambda_m$ , we can only represent the extremal values of  $\lambda_i$  exactly using the (4.6) ansatz whenever  $\lambda_1 \lambda_m < 0$ . Consequently, we have a slightly more complicated identity for general  $m > 2$

$$|\tilde{A}|_{\tilde{A}_0} = \frac{\max(0, \lambda_m) + \min(0, \lambda_1)}{\max(0, \lambda_m) - \min(0, \lambda_1)} \tilde{A} - \frac{2 \max(0, \lambda_m) \min(0, \lambda_1)}{\max(0, \lambda_m) - \min(0, \lambda_1)} \tilde{A}_0 - \sum_{i=2}^{m-1} f_i \tilde{r}_i \otimes \tilde{r}_i$$

with

$$\begin{aligned} f_i &= \frac{\max(0, \lambda_m) + \min(0, \lambda_1)}{\max(0, \lambda_m) - \min(0, \lambda_1)} \lambda_i - \frac{2 \max(0, \lambda_m) \min(0, \lambda_1)}{\max(0, \lambda_m) - \min(0, \lambda_1)} - |\lambda_i| \\ &= \frac{\max(0, \lambda_m) (\lambda_i^- - \min(0, \lambda_1)) + \min(0, \lambda_1) (\lambda_i^+ - \max(0, \lambda_m))}{\max(0, \lambda_m) - \min(0, \lambda_1)} \geq 0 . \end{aligned}$$

Next, consider the local path integral form of  $\mathbf{h}_{\text{SMV}}^d$

$$\begin{aligned} \mathbf{h}_{\text{SMV}}^d(\mathbf{v}_-, \mathbf{v}_+; \mathbf{n}) &= \int_0^1 |\tilde{A}(\bar{\mathbf{v}}(\theta))|_{\tilde{A}_0} [\mathbf{v}]_-^+ d\theta \\ &= \int_0^1 \left( \sigma_1(\theta) \tilde{A}(\bar{\mathbf{v}}(\theta); \mathbf{n}) - \sigma_2(\theta) \tilde{A}_0(\bar{\mathbf{v}}(\theta)) - \sum_{i=2}^{m-1} f_i(\theta) \tilde{r}_i \otimes \tilde{r}_i \right) [\mathbf{v}]_-^+ d\theta \end{aligned}$$

with

$$f_i(\theta) = \frac{\max(0, \lambda_m(\theta)) + \min(0, \lambda_1(\theta))}{\max(0, \lambda_m(\theta)) - \min(0, \lambda_1(\theta))} \lambda_i(\theta) - \frac{2 \max(0, \lambda_m(\theta)) \min(0, \lambda_1(\theta))}{\max(0, \lambda_m(\theta)) - \min(0, \lambda_1(\theta))} - |\lambda_i(\theta)| \geq 0$$

and

$$\begin{aligned} \sigma_1(\theta) &= \frac{\max(0, \lambda_m(\bar{\mathbf{v}}(\theta))) + \min(0, \lambda_1(\bar{\mathbf{v}}(\theta)))}{\max(0, \lambda_m(\bar{\mathbf{v}}(\theta))) - \min(0, \lambda_1(\bar{\mathbf{v}}(\theta)))} , \\ \sigma_2(\theta) &= \frac{2 \max(0, \lambda_m(\bar{\mathbf{v}}(\theta))) \min(0, \lambda_1(\bar{\mathbf{v}}(\theta)))}{\max(0, \lambda_m(\bar{\mathbf{v}}(\theta))) - \min(0, \lambda_1(\bar{\mathbf{v}}(\theta)))} . \end{aligned}$$

In addition, we will define the perturbed ratios<sup>2</sup>

$$\begin{aligned}\tilde{\sigma}_1(\theta) &= \frac{(\max(0, \lambda_m(\bar{\mathbf{v}}(\theta))) + \delta(\theta)) + (\min(0, \lambda_1(\bar{\mathbf{v}}(\theta))) - \gamma(\theta))}{(\max(0, \lambda_m(\bar{\mathbf{v}}(\theta))) + \delta(\theta)) - (\min(0, \lambda_1(\bar{\mathbf{v}}(\theta))) - \gamma(\theta))} , \\ \tilde{\sigma}_2(\theta) &= \frac{2(\max(0, \lambda_m(\bar{\mathbf{v}}(\theta))) + \delta(\theta)) (\min(0, \lambda_1(\bar{\mathbf{v}}(\theta))) - \gamma(\theta))}{(\max(0, \lambda_m(\bar{\mathbf{v}}(\theta))) + \delta(\theta)) - (\min(0, \lambda_1(\bar{\mathbf{v}}(\theta))) - \gamma(\theta))}\end{aligned}$$

for nonnegative bounded functions  $\delta(\theta) \geq 0$  and  $\gamma(\theta) \geq 0$ ,  $\theta \in [0, 1]$ . Examination of the scalar quantity

$$\begin{aligned}II &= \int_0^1 [\mathbf{v}]_-^+ \cdot \left( \tilde{\sigma}_1(\theta) \tilde{A}(\bar{\mathbf{v}}(\theta); \mathbf{n}) - \tilde{\sigma}_2(\theta) \tilde{A}_0(\bar{\mathbf{v}}(\theta)) \right) [\mathbf{v}]_-^+ d\theta \\ &\quad - \int_0^1 [\mathbf{v}]_-^+ \cdot \left( \sigma_1(\theta) \tilde{A}(\bar{\mathbf{v}}(\theta); \mathbf{n}) - \sigma_2(\theta) \tilde{A}_0(\bar{\mathbf{v}}(\theta)) \right) [\mathbf{v}]_-^+ d\theta \\ &= \int_0^1 [\mathbf{v}]_-^+ \tilde{R}(\theta) \cdot \left( (\tilde{\sigma}_1(\theta) - \sigma_1(\theta)) \Lambda(\theta) - (\tilde{\sigma}_2(\theta) - \sigma_2(\theta)) I_{m \times m} \right) \tilde{R}^T(\theta) [\mathbf{v}]_-^+ d\theta\end{aligned}$$

reveals that  $II \geq 0$  since for each component  $\lambda_i$  of  $\Lambda$ ,  $i = 1, \dots, m$ , (omitting the dependence on  $\theta$ )

$$\begin{aligned}(\tilde{\sigma}_1 - \sigma_1) \lambda_i - (\tilde{\sigma}_2 - \sigma_2) &= \frac{2 \gamma \max(0, \lambda_m) (\max(0, \lambda_m) + \delta - \lambda_i)}{((\max(0, \lambda_m) + \delta) - (\min(0, \lambda_1) - \gamma)) (\max(0, \lambda_m) - \min(0, \lambda_1))} \\ &\quad + \frac{2 \delta \min(0, \lambda_1) (\min(0, \lambda_1) - \gamma - \lambda_i)}{((\max(0, \lambda_m) + \delta) - (\min(0, \lambda_1) - \gamma)) (\max(0, \lambda_m) - \min(0, \lambda_1))} \\ &\geq 0 .\end{aligned}$$

Define infimum and supremum values of  $\min(0, \lambda_1(\theta))$  and  $\max(0, \lambda_m(\theta))$  respectively in the interval  $\theta \in [0, 1]$  as

$$\lambda_{\max} \equiv \sup_{0 \leq \theta \leq 1} \max(0, \lambda_m(\bar{\mathbf{v}}(\theta))) , \quad \lambda_{\min} \equiv \inf_{0 \leq \theta \leq 1} \min(0, \lambda_1(\bar{\mathbf{v}}(\theta)))$$

and set  $\epsilon(\theta)$  and  $\delta(\theta)$  as follows

$$\epsilon(\theta) = \lambda_{\max} - \max(0, \lambda_m(\bar{\mathbf{v}}(\theta))) \geq 0 , \quad \delta(\theta) = \min(0, \lambda_1(\bar{\mathbf{v}}(\theta))) - \lambda_{\min} \geq 0 .$$

This renders  $\tilde{\sigma}_1(\theta)$  and  $\tilde{\sigma}_2(\theta)$  now  $\theta$ -independent, i.e.

$$\tilde{\sigma}_1 = \frac{\lambda_{\max} + \lambda_{\min}}{\lambda_{\max} - \lambda_{\min}} , \quad \tilde{\sigma}_2 = \frac{2\lambda_{\max}\lambda_{\min}}{\lambda_{\max} - \lambda_{\min}} .$$

Consequently, for  $II \geq 0$  and  $f_i(\theta) \geq 0$  we have

$$\begin{aligned}[\mathbf{v}]_-^+ \cdot \mathbf{h}_{\text{SMV}}^d &= [\mathbf{v}]_-^+ \cdot \int_0^1 \left( \sigma_1(\theta) \tilde{A}(\bar{\mathbf{v}}(\theta); \mathbf{n}) - \sigma_2(\theta) \tilde{A}_0(\bar{\mathbf{v}}(\theta)) \right) [\mathbf{v}]_-^+ d\theta \\ &\quad - [\mathbf{v}]_-^+ \cdot \int_0^1 \sum_{i=2}^{m-1} f_i(\theta) \tilde{\mathbf{r}}_i \otimes \tilde{\mathbf{r}}_i [\mathbf{v}]_-^+ d\theta \\ &\leq [\mathbf{v}]_-^+ \cdot \int_0^1 \left( \tilde{\sigma}_1(\theta) \tilde{A}(\bar{\mathbf{v}}(\theta); \mathbf{n}) - \tilde{\sigma}_2(\theta) \tilde{A}_0(\bar{\mathbf{v}}(\theta)) \right) [\mathbf{v}]_-^+ d\theta \\ &\quad - [\mathbf{v}]_-^+ \cdot \int_0^1 \sum_{i=2}^{m-1} f_i(\theta) \tilde{\mathbf{r}}_i \otimes \tilde{\mathbf{r}}_i [\mathbf{v}]_-^+ d\theta \\ &= [\mathbf{v}]_-^+ \cdot \left( \tilde{\sigma}_1 \int_0^1 \tilde{A}(\bar{\mathbf{v}}(\theta); \mathbf{n}) [\mathbf{v}]_-^+ d\theta - \tilde{\sigma}_2 \int_0^1 \tilde{A}_0(\bar{\mathbf{v}}(\theta)) [\mathbf{v}]_-^+ d\theta \right)\end{aligned}$$

<sup>2</sup>Our strategy will be to later define the nonnegative function  $\delta(\theta)$  as the “gap” between local value of  $\max(0, \lambda_m(\theta))$  and the supremum value in the interval  $\theta \in [0, 1]$  and similarly  $\gamma(\theta)$  will represent the gap between the  $\min(0, \lambda_1(\theta))$  and the infimum value.

$$\begin{aligned}
& -[\mathbf{v}]_{-}^{+} \cdot \int_0^1 \sum_{i=2}^{m-1} f_i(\theta) \tilde{\mathbf{r}}_i \otimes \tilde{\mathbf{r}}_i [\mathbf{v}]_{-}^{+} d\theta \\
& = [\mathbf{v}]_{-}^{+} \cdot \left( \tilde{\sigma}_1 [\mathbf{f}(\mathbf{v}; \mathbf{n})]_{-}^{+} - [\mathbf{v}]_{-}^{+} \cdot \tilde{\sigma}_2 [\mathbf{u}(\mathbf{v})]_{-}^{+} \right) \\
& \quad - [\mathbf{v}]_{-}^{+} \cdot \int_0^1 \sum_{i=2}^{m-1} f_i(\theta) \tilde{\mathbf{r}}_i \otimes \tilde{\mathbf{r}}_i [\mathbf{v}]_{-}^{+} d\theta \\
& = [\mathbf{v}]_{-}^{+} \cdot \mathbf{h}_{\text{SHLLEM}}^d \\
& \leq [\mathbf{v}]_{-}^{+} \cdot \left( \tilde{\sigma}_1 [\mathbf{f}(\mathbf{v}; \mathbf{n})]_{-}^{+} - [\mathbf{v}]_{-}^{+} \cdot \tilde{\sigma}_2 [\mathbf{u}(\mathbf{v})]_{-}^{+} \right) \\
& = [\mathbf{v}]_{-}^{+} \cdot \mathbf{h}_{\text{SHLLE}}^d .
\end{aligned}$$

which completes the proof for the SHLLE and SHLLEM fluxes. ■

Discrete Kinetic Symmetric Mean-Value Flux (DKSMV):

$$\begin{aligned}
[\mathbf{v}]_{-}^{+} \cdot \mathbf{h}_{\text{KSMV}}^d & = [\mathbf{v}]_{-}^{+} \cdot \int_0^1 \langle |v \cdot \mathbf{n}| \mathbf{m} \otimes \mathbf{m} \exp(\bar{\mathbf{v}}(\theta) \cdot \mathbf{m}(v)) \rangle [\mathbf{v}]_{-}^{+} d\theta \\
& = [\mathbf{v}]_{-}^{+} \cdot \left\langle |v \cdot \mathbf{n}| \mathbf{m} \otimes \mathbf{m} \int_0^1 \exp(\bar{\mathbf{v}}(\theta) \cdot \mathbf{m}(v)) d\theta \right\rangle [\mathbf{v}]_{-}^{+} \\
& = \left\langle |v \cdot \mathbf{n}| ([\mathbf{v}]_{-}^{+} \cdot \mathbf{m})^2 \int_0^1 \exp(\bar{\mathbf{v}}(\theta) \cdot \mathbf{m}(v)) d\theta \right\rangle .
\end{aligned}$$

Considering the scalar function

$$g(\theta) = \exp(\bar{\mathbf{v}}(\theta) \cdot \mathbf{m}(v))$$

followed by twice differentiation

$$g''(\theta) = ([\mathbf{v}]_{-}^{+} \cdot \mathbf{m}(v))^2 \exp(\bar{\mathbf{v}}(\theta) \cdot \mathbf{m}(v)) \geq 0 .$$

Hence,  $g(\theta)$  is yet another convex function so that

$$g(\theta) \leq (1 - \theta) g(0) + \theta g(1)$$

and finally

$$\begin{aligned}
[\mathbf{v}]_{-}^{+} \cdot \mathbf{h}_{\text{KSMV}}^d & = [\mathbf{v}]_{-}^{+} \cdot \int_0^1 \langle |v \cdot \mathbf{n}| \mathbf{m} \otimes \mathbf{m} \exp(\bar{\mathbf{v}}(\theta) \cdot \mathbf{m}(v)) \rangle [\mathbf{v}]_{-}^{+} d\theta \\
& \leq [\mathbf{v}]_{-}^{+} \cdot \frac{1}{2} \langle |v \cdot \mathbf{n}| \mathbf{m} \otimes \mathbf{m} (\exp(\mathbf{v}_{-} \cdot \mathbf{m}(v)) + \exp(\mathbf{v}_{+} \cdot \mathbf{m}(v))) \rangle [\mathbf{v}]_{-}^{+} \\
& = [\mathbf{v}]_{-}^{+} \cdot \mathbf{h}_{\text{DKSMV}}^d \quad \blacksquare
\end{aligned}$$

**5. Concluding Remarks.** The analysis of this paper confirms energy stability for several numerical flux functions that are of practical merit when used in computational fluid dynamics computations. Even so, the theoretical framework developed here applies more generally and has application to many nonlinear conservation law systems with entropy extensions that are not explicitly discussed here. In these settings, the analysis presented may be invaluable because there may not be the large body of numerical methods developed before hand to guide the development of new numerical fluxes, discretization, and stabilization. Our general goal is to pursue these new problem areas in forthcoming work.

**Appendix A. The  $\exp(\mathbf{v} \cdot \mathbf{m}(v, I))$  Boltzmann Moment Structure for a  $\gamma$ -Law Equation of State.** For a  $\gamma$ -law (polytropic) gas, one has

$$p = (\gamma - 1)\rho\epsilon, \quad T = (\gamma - 1)\epsilon .$$



Following Perthame [18], we consider the following Maxwellian in  $\mathbb{R}^d$  for a  $\gamma$ -law gas:

$$(A.1) \quad f(\rho, \mathbf{u}, T; \mathbf{v}, I) = \frac{\rho}{\alpha(\gamma, d) T^{d/2+1/\delta}} e^{-(|\mathbf{u}-\mathbf{v}|^2/2+I^\delta)/T}$$

with

$$\delta = \frac{1}{\frac{1}{\gamma-1} - \frac{d}{2}}$$

and

$$\alpha(\gamma, d) = \int_{\mathbb{R}^d} e^{-|\mathbf{v}|^2/2} d\mathbf{v} \cdot \int_{\mathbb{R}^+} e^{-I^\delta} dI .$$

Using this particular form, Perthame shows that the Euler equations for a  $\gamma$ -law gas are obtained as the following moments

$$(A.2) \quad \mathbf{m}(\mathbf{v}, I) = \begin{pmatrix} 1 \\ \mathbf{v} \\ |\mathbf{v}|^2/2 + I^\delta \end{pmatrix}$$

$$\mathbf{u} = \langle \mathbf{m} f \rangle, \quad \mathbf{f}^i = \langle v_i \mathbf{m} f \rangle, \quad \langle \cdot \rangle \equiv \int_{\mathbb{R}^d} \int_{\mathbb{R}^+} (\cdot) dI d\mathbf{v} .$$

The nonobvious energy moment  $|\mathbf{v}|^2/2 + I^\delta$  was devised by Perthame rather than the more standard moment  $|\mathbf{v}|^2/2 + I$  (see for example [6, 7]) in order that a classical Boltzmann entropy  $H(f) = \int f \log f$  be obtained.

Let us now verify that this choice of moments yields an exponential form for the conjugate entropy function of the form

$$\mathcal{U}(\mathbf{v}) = \langle f \rangle = \langle c(\gamma, d) \exp(\mathbf{v} \cdot \mathbf{m}(\mathbf{v}, I)) \rangle, \quad c(\gamma, d) > 0$$

which is sufficient for our purposes. Inserting the expression for  $\delta$  into the temperature term appearing in the Maxwellian yields

$$(A.3) \quad f(\rho, \mathbf{u}, T; \mathbf{v}, I) = \frac{1}{\alpha(\gamma, d)} \frac{\rho}{T^{1/(\gamma-1)}} e^{-(|\mathbf{u}-\mathbf{v}|^2/2+I^\delta)/T}$$

and compare this with the expression  $\exp(\mathbf{v} \cdot \mathbf{m})$  obtained using the entropy function<sup>3</sup>

$$U(\mathbf{u}) = -\frac{\rho s}{(\gamma-1)}$$

so that

$$\mathbf{v} = U_{,\mathbf{u}}^T = \begin{pmatrix} -\frac{s}{\gamma-1} + \frac{\gamma}{\gamma-1} - \frac{|\mathbf{u}|^2}{2T} \\ \frac{\mathbf{u}}{T} \\ -\frac{1}{T} \end{pmatrix} = \begin{pmatrix} \log\left(\frac{\rho}{T^{1/(\gamma-1)}}\right) + \frac{\gamma}{\gamma-1} - \frac{|\mathbf{u}|^2}{2T} \\ \frac{\mathbf{u}}{T} \\ -\frac{1}{T} \end{pmatrix}$$

and finally

$$\exp(\mathbf{v} \cdot \mathbf{m}(\mathbf{v}, I)) = e^{-\gamma/(\gamma-1)} \frac{\rho}{T^{1/(\gamma-1)}} e^{-(|\mathbf{u}-\mathbf{v}|^2/2+I^\delta)/T} .$$

Comparing with (A.1), we obtain the exponential form for the conjugate entropy function

$$\mathcal{U}(\mathbf{v}) = \langle f \rangle = \langle c(\gamma, d) \exp(\mathbf{v} \cdot \mathbf{m}(\mathbf{v}, I)) \rangle$$

with

$$c(\gamma, d) = \frac{e^{\gamma/(\gamma-1)}}{\alpha(\gamma, d)} > 0 .$$

<sup>3</sup>The choice of  $1/(\gamma-1)$  scaling of the entropy function comes from our desire to match the  $\rho/T^{1/(\gamma-1)}$  term appearing in Perthame's Maxwellian